



ELSEVIER

Available online at www.sciencedirect.com

Journal of Computational and Applied Mathematics 199 (2007) 418–423

JOURNAL OF
COMPUTATIONAL AND
APPLIED MATHEMATICSwww.elsevier.com/locate/cam

Interval versions of statistical techniques with applications to environmental analysis, bioinformatics, and privacy in statistical databases

Vladik Kreinovich^{a,*}, Luc Longpré^a, Scott A. Starks^a, Gang Xiang^a, Jan Beck^a,
Raj Kandathi^a, Asis Nayak^a, Scott Ferson^b, Janos Hajagos^{b,c}

^aNASA Pan-American Center for Earth and Environmental Studies (PACES), University of Texas, El Paso, TX 79968, USA

^bApplied Biomathematics, 100 North Country Road, Setauket, NY 11733, USA

^cDepartment of Ecology and Evolution, State University of New York, Stony Brook, NY 11794, USA

Received 27 December 2004

Abstract

In many areas of science and engineering, it is desirable to estimate statistical characteristics (mean, variance, covariance, etc.) under interval uncertainty. For example, we may want to use the measured values $x(t)$ of a pollution level in a lake at different moments of time to estimate the average pollution level; however, we do not know the exact values $x(t)$ —e.g., if one of the measurement results is 0, this simply means that the actual (unknown) value of $x(t)$ can be anywhere between 0 and the detection limit (DL). We must, therefore, modify the existing statistical algorithms to process such interval data.

Such a modification is also necessary to process data from statistical databases, where, in order to maintain privacy, we only keep interval ranges instead of the actual numeric data (e.g., a salary range instead of the actual salary).

Most resulting computational problems are NP-hard—which means, crudely speaking, that in general, no computationally efficient algorithm can solve all particular cases of the corresponding problem. In this paper, we overview practical situations in which computationally efficient algorithms exist: e.g., situations when measurements are very accurate, or when all the measurements are done with one (or few) instruments.

As a case study, we consider a practical problem from bioinformatics: to discover the genetic difference between the cancer cells and the healthy cells, we must process the measurements results and find the concentrations c and h of a given gene in cancer and in healthy cells. This is a particular case of a general situation in which, to estimate states or parameters which are not directly accessible by measurements, we must solve a system of equations in which coefficients are only known with interval uncertainty. We show that in general, this problem is NP-hard, and we describe new efficient algorithms for solving this problem in practically important situations.

© 2006 Elsevier B.V. All rights reserved.

Keywords: Intervals and probabilities; Environmental analysis; Bioinformatics; Privacy; Statistical databases

* Corresponding author. Tel.: +1 915 747 6951; fax: +1 915 747 5030.

E-mail address: vladik@utep.edu (V. Kreinovich).

1. Statistical analysis is important

Many aspects of engineering and science involve statistical uncertainty. It is, therefore, desirable to estimate statistical characteristics such as mean, variance, covariance, etc., i.e., compute statistics such as $E(x) = (1/n)(x_1 + \dots + x_n)$, $V(x) = (1/(n-1)) \cdot \sum_{i=1}^n (x_i - E(x))^2$, and $C(x, y) = (1/(n-1)) \cdot \sum_{i=1}^n (x_i - E(x)) \cdot (y_i - E(y))$. For example, in *non-destructive testing*, outliers are indications of faults; outliers are often detected as values outside the interval $[E(x) - k_0 \cdot \sqrt{V(x)}, E(x) + k_0 \cdot \sqrt{V(x)}]$ for $k_0 = 2, 3$, or 6 . In *geophysics*, outliers indicate possible locations of minerals. In *biomedical systems*, statistical analysis often leads to improvements in medical recommendations.

Comment: In many practical situations, e.g., when measuring the magnitude and orientation of a magnetic field \vec{H} , what we measure is not a single-component (*scalar*) value $x \in R$, but a *multi-component* value: e.g., a vector $\vec{H} \in R^3$. In such situations, it is reasonable to estimate, e.g., the mean value of the corresponding vector measurements as $E(\vec{H}) = (1/n) \cdot (\vec{H}_1 + \dots + \vec{H}_n)$.

From the physical viewpoint, statistical analysis of the vector data is different from the statistical analysis of the scalar data. However, from the purely computational viewpoint, the problem is largely the same: e.g., for each coordinate α , the α -component $E_\alpha(\vec{H})$ of the average vector $E(\vec{H})$ is equal to the arithmetic average of the corresponding components of \vec{H}_i . Since our objective is to help in computations, in the following text, we will limit our description to scalar values $x_i \in R$.

2. Interval uncertainty

Traditional statistics assumes that we know the exact sample values x_1, \dots, x_n . In practice, often, we only know x_i with interval uncertainty: $x_i \in [\underline{x}_i, \bar{x}_i]$ see, e.g., [2].

For example, values x_i usually come from measurements, and we often only know the upper bounds Δ_i on the measurement error $\Delta x_i \stackrel{\text{def}}{=} \tilde{x}_i - x_i$. So, the only information that we have about x_i is that $x_i \in [\tilde{x}_i - \Delta_i, \tilde{x}_i + \Delta_i]$.

Another source of interval uncertainty is the existence of detection limits for different sensors: if a sensor, e.g., did not detect any ozone, this means that the ozone concentration is below its detection limit (DL), i.e., in the interval $[0, \text{DL}]$.

Yet another source of interval uncertainty is discretized data: if we experiment on the fish and watch it daily, and a fish is alive on Day 5 but dead on Day 6, then all we know about its lifetime is that it is in the interval $[5, 6]$.

Expert estimates often come as intervals.

The need to keep privacy in statistical (e.g., medical) databases also often leads to the fact that instead of recording, e.g., exact age, what we only record is the interval $[40, 50]$.

Summarizing, often, instead of the actual values x_1, \dots, x_n , we only know the intervals $\mathbf{x}_1 = [\underline{x}_1, \bar{x}_1], \dots, \mathbf{x}_n = [\underline{x}_n, \bar{x}_n]$ that contain x_i . Different values $x_i \in \mathbf{x}_i$ lead to different values of the statistic $S(x_1, \dots, x_n)$. It is desirable to find the range of such values:

$$S(\mathbf{x}_1, \dots, \mathbf{x}_n) \stackrel{\text{def}}{=} \{S(x_1, \dots, x_n) | x_1 \in \mathbf{x}_1, \dots, x_n \in \mathbf{x}_n\}.$$

3. Simple and hard cases

The mean $E(x)$ is monotonic, so $\mathbf{E}(x) = [\underline{E}(x), \bar{E}(x)]$, where $\underline{E}(x) = (1/n)(\underline{x}_1 + \dots + \underline{x}_n)$ and $\bar{E}(x) = (1/n)(\bar{x}_1 + \dots + \bar{x}_n)$.

For other statistics such as variance $V(x)$ or covariance $C(x, y)$, the problem is, in general, NP-hard [1,3,5]. In such cases, in general, we have to use approximate techniques.

4. Linearization and its limitations

One of the known approximate techniques is linearization, when we approximate the statistics S with the linear terms in its Taylor expansion: $S \approx S_{\text{lin}} = S_0 - \sum_{i=1}^n S_i \cdot \Delta x_i$, where $S_0 \stackrel{\text{def}}{=} S(\tilde{x}_1, \dots, \tilde{x}_n)$, $S_i \stackrel{\text{def}}{=} (\partial S / \partial x_i)(\tilde{x}_1, \dots, \tilde{x}_n)$, and $\Delta x_i \stackrel{\text{def}}{=} \tilde{x}_i - x_i$. For the linear function, we get the exact formula for the range: $\mathbf{S} = [S_0 - \Delta_S, S_0 + \Delta_S]$, where $\Delta_S \stackrel{\text{def}}{=} \sum_{i=1}^n |S_i| \cdot \Delta_i$.

However, linearization is not always acceptable. Sometimes, the intervals are wide, so that quadratic terms cannot be ignored. Sometimes—e.g., in cases of bioregulations—we want to *guarantee* that, e.g., the variance $V(x)$ is below a given threshold V_0 . So, we need validated techniques.

Since we cannot provide efficient algorithms for the general case, we must find practically useful cases for which an efficient algorithm is possible.

5. Classes of problems for which efficient algorithms are known:

1. *Narrow intervals*: no two intervals x_i intersect.
2. *Slightly wider intervals*: for some integer K , no set of K intervals has a common intersection.
3. *Single measuring instrument (MI)*: no two intervals are subsets of each other, i.e., $[\underline{x}_i, \bar{x}_i] \not\subseteq [\underline{x}_j, \bar{x}_j]$ (non-degenerate results are allowed).
4. *Same accuracy measurement*: $\Delta_1 = \dots = \Delta_n$.
5. *Several MI*: intervals are divided into several subgroups each of which comes from a single MI.
6. *Privacy case*: intervals are formed from the given partition, e.g., 10–20, 20–30, etc.; in this case, every two non-degenerate intervals either coincide or do not intersect.
7. *Non-detects*: every measurement result is either an exact value or a *non-detect*, i.e., an interval $[0, DL_i]$ for some real number DL_i .

In these cases, we have the following complexity results [4,6], where Class 0 means the general case (when almost all problems are NP-hard),

$$L(x) \stackrel{\text{def}}{=} E(x) - k_0 \cdot \sqrt{V(x)}, \quad U(x) \stackrel{\text{def}}{=} E(x) + k_0 \cdot \sqrt{V(x)},$$

$R(x)$ is the largest value k_0 for which $x_0 \notin [L(x), U(x)]$, where x_0 is a given value, i.e., $R(x) \stackrel{\text{def}}{=} |x_0 - E(x)|/\sqrt{V(x)}$, and $M_m(x)$ is m th central moment: $M_m(x) \stackrel{\text{def}}{=} (1/n) \sum_{i=1}^n |x_i - E(x)|^m$.

#	$E(x)$	$V(x), L(x), U(x), R(x), M_{2p}(x)$	$C(x, y)$	$M_{2p+1}(x)$
0	$O(n)$	NP-hard	NP-hard	?
1	$O(n)$	$O(n \cdot \log(n))$	$O(n^2)$	$O(n^2)$
2	$O(n)$	$O(n \cdot \log(n))$	$O(n^2)$	$O(n^2)$
3	$O(n)$	$O(n \cdot \log(n))$?	?
4	$O(n)$	$O(n \cdot \log(n))$	$O(n^3)$?
5	$O(n)$	$O(n^m)$?	?
6	$O(n)$	$O(n \cdot \log(n))$	$O(n^2)$?
7	$O(n)$	$O(n \cdot \log(n))$?	?

6. Case when only d out of n data points are non-degenerate intervals

In this case, we have the following complexity results:

#	$E(x)$	$V(x), L(x), U(x), R(x), M_{2p}(x)$	$C(x, y)$	$M_{2p+1}(x)$
0	$O(n)$	NP-hard	NP-hard	?
1	$O(n)$	$O(n + d \cdot \log(d))$	$O(n + d^2)$	$O(n + d^2)$
2	$O(n)$	$O(n + d \cdot \log(d))$	$O(n + d^2)$	$O(n + d^2)$
3	$O(n)$	$O(n + d \cdot \log(d))$?	?
4	$O(n)$	$O(n + d \cdot \log(d))$	$O(n + d^3)$?
5	$O(n)$	$O(n + d^m)$?	?
6	$O(n)$	$O(n + d \cdot \log(d))$	$O(n + d^2)$?
7	$O(n)$	$O(n + d \cdot \log(d))$?	?

7. Other statistics

We have mentioned that an important source of interval uncertainty is the existence of the lower detection limits for sensors: if a sensor does not detect any signal this means that the actual value of the measured quantity is below its DL, i.e., in the interval $[0, \text{DL}]$.

Another practically important source of uncertainty is the fact that many sensors also have saturation values x_{\max} : if the sensor registers the value $\tilde{x}_i = x_{\max}$, then the only information that we know about the true value x is that $x \geq x_{\max}$, i.e., that $x \in [x_{\max}, \infty)$. If one of the measurements \tilde{x}_i is equal to the saturation value, then, e.g., the arithmetic average $E(x) = (1/n) \cdot (x_1 + \dots + x_n)$ of the actual values x_i can be arbitrarily large.

For such situations, we need to use different methods for estimating the expected value (mean) $E\{x\}$ of a random variable from the sample x_1, \dots, x_n . One such method is a median. Median is a particular case of an important class of statistical *L-estimates*: we order the values x_i into a (non-strictly) increasing sequence $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$, and then estimate $E\{x\}$ as $\sum_{i=1}^n w_i \cdot x_{(i)}$.

Alternative methods for estimating $E\{x\}$ are also useful in other practical situations—e.g., if, in addition to measurement results, the values x_i contain erroneously recorded values. Other widely used alternative methods for estimating $E\{x\}$ include [7,8]:

- *weighted mean* that is defined by the condition $\sum_{i=1}^n (x_i - E)^2 / \sigma^2 \rightarrow \min_E$, so

$$E_w = \sum_{i=1}^n p_i \cdot x_i, \quad \text{where } p_i \stackrel{\text{def}}{=} \frac{\sigma_i^{-2}}{\sum_{j=1}^n \sigma_j^{-2}};$$

- *M-estimates*: $\sum_{i=1}^n \psi(|x_i - a|) \rightarrow \max_a$ for some function $\psi(x)$; average is a particular case of an M-estimate, corresponding to $\psi(x) = x^2$.

They are all monotonic functions of x_i , so their ranges can be computed in time $O(n)$.

8. Case study: bioinformatics

In cancer research, it is important to find out the genetic difference between the cancer cells and the healthy cells. In the ideal world, we should be able to have a sample of cancer cells, and a sample of healthy cells, and thus directly measure the concentrations c and h of a given gene in cancer and in healthy cells. In reality, it is very difficult to separate the cells, so we have to deal with samples that contain both cancer and normal cells. Let y_i denote the result of measuring the concentration of the gene in i th sample, and let x_i denote the percentage of cancer cells in i th sample. Then, we should have $x_i \cdot c + (1 - x_i) \cdot h \approx y_i$ [9] (approximately equal because there are measurement errors in measuring y_i).

Let us first consider an idealized case in which we know the exact percentages x_i . In this case, we can find the desired values c and h by solving a system of linear equations $x_i \cdot c + (1 - x_i) \cdot h \approx y_i$ with two unknowns c and h .

It is worth mentioning that this system can be somewhat simplified if instead of c , we consider a new variable $a \stackrel{\text{def}}{=} c - h$. In terms of the new unknowns a and h , the system takes the following form: $a \cdot x_i + h \approx y_i$.

The errors of measuring y_i are normally i.i.d. random variables, so to estimate a and h , we can use the least squares method (LSM) $\sum_{i=1}^n (a \cdot x_i + h - y_i)^2 \rightarrow \min_{a,h}$, according to which $a = C(x, y) / V(x)$ and $h = E(y) - a \cdot E(x)$. Once we know $a = c - h$ and h , we can then estimate c as $a + h$.

The problem is that the concentrations x_i come from experts who manually count different cells, and experts can only provide interval bounds on the values x_i such as $x_i \in [0.7, 0.8]$. Different values of x_i in the corresponding intervals lead to different values of a and h . It is, therefore, desirable to find the range of a and h corresponding to all possible values $x_i \in [\underline{x}_i, \bar{x}_i]$.

Comment: Our motivation for solving this problem comes from bioinformatics, but similar problems appear in various practical situations where measurements with uncertainties are available and statistical data are to be processed.

9. Linear approximation

Let $\tilde{x}_i = (\underline{x}_i + \bar{x}_i)/2$ be the midpoint of i th intervals, and let $\Delta_i = (\bar{x}_i - \underline{x}_i)/2$ be its half-width. For a , we have

$$\frac{\partial a}{\partial x_i} = \frac{1}{(n-1) \cdot V(x)} \cdot (y_i - E(y) - 2a \cdot x_i + 2a \cdot E(x)).$$

We can use the formula $E(y) = a \cdot E(x) + h$ to simplify this expression, resulting in

$$\Delta_a = (1/((n-1) \cdot V(x))) \sum_{i=1}^n |\Delta y_i - a \cdot \Delta x_i| \cdot \Delta_i,$$

where we denoted $\Delta y_i \stackrel{\text{def}}{=} y_i - a \cdot x_i - h$ and $\Delta x_i \stackrel{\text{def}}{=} x_i - E(x)$.

Since $h = E(y) - a \cdot E(x)$, we have $\partial h / \partial x_i = -\partial a / \partial x_i \cdot E(x) - 1/n \cdot a$, so $\Delta_h = \sum_{i=1}^n |\partial h / \partial x_i| \cdot \Delta_i$.

10. Prior estimation of the resulting accuracy

The above formulas provide us with the accuracy *after* the data have been processed. It is often desirable to have an estimate *prior* to measurements, to make sure that we will get c and h with desired accuracy.

The difference Δy_i is a measurement error, so it is normally distributed with 0 mean and standard deviation $\sigma(y)$ corresponding to the accuracy of measuring y_i . The difference Δx_i is distributed with 0 mean and standard deviation $\sqrt{V(x)}$. For estimation purposes, it is reasonable to assume that the values Δx_i are also normally distributed. It is also reasonable to assume that the errors in x_i and y_i are uncorrelated, so the linear combination $\Delta y_i - a \cdot \Delta x_i$ is also normally distributed, with 0 mean and variance $\sigma_y^2 + a^2 \cdot V(x)$. It is also reasonable to assume that all the values Δ_i are approximately the same: $\Delta_i \approx \Delta$.

For normal distribution ξ with 0 mean and standard deviation σ , the mean value of $|\xi|$ is equal to $\sqrt{2/\pi} \cdot \sigma$. Thus, the absolute value $|\Delta y_i - a \cdot \Delta x_i|$ of the above combination has a mean value $\sqrt{2/\pi} \cdot \sqrt{\sigma_y^2 + a^2 \cdot V(x)}$. Hence, the expected value of Δ_a is equal to $(2/\pi) \cdot \sqrt{\sigma_y^2 + a^2 \cdot V(x)} \cdot \Delta / V(x)$.

Since measurements are usually more accurate than expert estimates, we have $\sigma_y^2 \ll V(x)$, hence $\Delta_a \approx (2/\pi) \cdot a \cdot \Delta$. Similar estimates can be given for Δ_h .

11. In general, finding the exact range is NP-hard

Let us show that in general, finding the exact range for the ratio $C(x, y)/V(x)$ is an NP-hard problem.

The proof is similar to the proof that computing the range for the variance is NP-hard [1,3,5]: namely, we reduce a partition problem (known to be NP-hard) to our problem. In the partition problem, we are given m positive integers s_1, \dots, s_m , and we must check whether there exist values $\varepsilon_i \in \{-1, 1\}$ for which $\sum_{i=1}^m \varepsilon_i \cdot s_i = 0$. We will reduce this problem to the following problem: $n = m + 2$, $y_1 = \dots = y_m = 0$, $y_{m+1} = 1$, $y_{m+2} = -1$, $x_i = [-s_i, s_i]$ for $i \leq m$, $x_{m+1} = 1$, and $x_{m+2} = -1$. In this case, $E(y) = 0$, so

$$C(x, y) = \frac{1}{n-1} \sum_{i=1}^n x_i \cdot y_i - \frac{n}{n-1} \cdot E(x) \cdot E(y) = \frac{2}{m+2}.$$

Therefore, $C(x, y)/V(x) \rightarrow \min$ if and only if $V(x) \rightarrow \max$.

Here,

$$V(x) = \frac{1}{m+1} \cdot \left(\sum_{i=1}^m x_i^2 + 2 \right) - \frac{m+2}{m+1} \cdot \left(\frac{1}{m+2} \cdot \sum_{i=1}^m x_i \right)^2.$$

Since $|x_i| \leq s_i$, we always have $V(x) \leq V_0 \stackrel{\text{def}}{=} (1/(m+1)) \cdot (\sum_{i=1}^m s_i^2 + 2)$, and the only possibility to have $V(x) = V_0$ is when $x_i = \pm s_i$ for all i and $\sum x_i = 0$. Thus, $V(x) = V_0$ if and only if the original partition problem has a solution. Hence, $C(x, y)/V(x) = 2/(\sum s_i^2 + 2)$ if and only if the original instance of the partition problem has a solution.

The reduction is proven, so our problem is indeed NP-hard.

Comment: In this proof, we consider the case when the values x_i can be negative and larger than 1, while in bioinformatics, x_i is always between 0 and 1. However, we can easily modify this proof: first, we can shift all the values x_i by the same constant to make them positive; shift does not change neither $C(x, y)$ nor $V(x)$. Second, to make the positive values ≤ 1 , we can then re-scale the values x_i ($x_i \rightarrow \lambda \cdot x_i$), thus multiplying $C(x, y)/V(x)$ by a known constant.

As a result, we get new values $x'_i = \frac{1}{2} \cdot (1 + x_i/K)$, where $K \stackrel{\text{def}}{=} \max s_i$, for which $x'_i \in [0, 1]$ and the problem of computing $C(x, y)/V(x)$ is still NP-hard.

12. What can we do?

One possibility is to use known algorithms to find the ranges for $C(x, y)$ and for $V(x)$, and then use the division operation from interval arithmetic to get the interval that is guaranteed to contain $C(x, y)/V(x)$.

Acknowledgments

This work was supported by NASA under cooperative agreement NCC5-209, by NSF Grants EAR-0112968, EAR-0225670, and EIA-0321328, by Army Research Laboratories Grant DATM-05-02-C-0046, and by NIH Grant 3T34GM008048-20S1.

The authors are very thankful to Ilya Shmulevich from the University of Texas M. D. Anderson Cancer Center for the formulation of the case study problem and valuable discussions, and to the anonymous referees for useful suggestions.

References

- [1] S. Ferson, L. Ginzburg, V. Kreinovich, L. Longpré, M. Aviles, Exact bounds on finite populations of interval data, *Reliab. Comput.* 11 (3) (2005) 207–233.
- [2] L. Jaulin, M. Keiffer, O. Didrit, E. Walter, *Applied Interval Analysis*, Springer, Berlin, 2001.
- [3] V. Kreinovich, Probabilities, intervals, what next? Optimization problems related to extension of interval computations to situations with partial information about probabilities, *J. Global Optim.* 29 (3) (2004) 265–280.
- [4] V. Kreinovich, J. Beck, C. Ferregut, A. Sanchez, G.R. Keller, M. Averill, S.A. Starks, Monte-Carlo-type techniques for processing interval uncertainty, and their engineering applications, *Proceedings of the Workshop on Reliable Engineering Computing*, Savannah, GA, September 15–17, 2004, pp. 139–160.
- [5] V. Kreinovich, L. Longpré, Computational complexity and feasibility of data processing and interval computations, with extension to cases when we have partial information about probabilities, in: V. Brattka, M. Schröder, K. Weihrauch, N. Zhong (Eds.), *Proceedings of Conference on Computability and Complexity in Analysis CCA'2003*, Cincinnati, OH, USA, August 28–30, 2003, pp. 19–54.
- [6] V. Kreinovich, G. Xiang, S.A. Starks, L. Longpré, M. Ceberio, R. Araiza, J. Beck, R. Kandathi, A. Nayak, R. Torres, J. Hajagos, Towards combining probabilistic and interval uncertainty in engineering calculations: algorithms for computing statistics under interval uncertainty, and their computational complexity, *Reliab. Comput.*, to appear.
- [7] A. Papoulis, *Probability, Random Variables, and Stochastic Processes*, McGraw-Hill, New York, 2002.
- [8] H.M. Wadsworth Jr., *Handbook of Statistical Methods for Engineers and Scientists*, McGraw-Hill, New York, 1990.
- [9] W. Zhang, I. Shmulevich, J. Astola, *Microarray Quality Control*, Wiley, Hoboken, NJ, 2004.